UCLA VISIONLAB

Visual-Inertial-Semantic Scene Representation for 3D Object Detection

Introduction

Overview

- We describe a system to detect objects in three-dimensional space using visual and inertial sensors (accelerometer and gyroscope).
- The resulting system can process the video stream causally in real time, and provides a representation of objects in the scene that is persistent.

Motivation

Principles

- Objects exist in the scene, not in the image;
- They persist, so confidence on their presence should grow as more evidence is accrued from multiple (test) images;
- Once seen, the system should be aware of their presence even when temporarily not visible;
- Such awareness should allow it to predict when they will return into view, based on scene geometry and topology;
- Objects have characteristic shape and size in 3D, and vestibular (*inertial*) sensors provide a global scale and orientation reference that the system should leverage on.

Methods

Problem Formalization

Given measurements up to time t y^t Estimate scene ξ and objects z^{j} with geometric (shape & pose) s_j and semantic (class) l_j attributes Quantity of interest: posterior of objects in the scene $p(\xi, z^j | y^t) = p(\xi | z^j) p(z^j | y^t)$ Context Learned from data

which is a *minimal sufficient* representation.

Methods (cont'd)



- Semantics: Nvidia GTX 760 Overall ~17 FPS
- Bottleneck: image-based object detectors
- Current implementation runs CNN every 3 frames

Jingming Dong, Xiaohan Fei, Stefano Soatto

	Comp	bar	risor	n ar	d	Eval	uat	io	n	
	Position error	< 0.5 m			< 1 m			$< 1.5 { m m}$		
Orientation error	method	#TP	Precision	Recall	#TP	Precision	Recall	#TP	Precision	Recall
< 30°	Ours-FNL	150	0.14	0.10	355	0.34	0.24	513	0.49	0.35
	Ours-INST	135	0.13	0.09	270	0.26	0.18	368	0.35	0.25
	SubCNN	99	0.10	0.07	254	0.26	0.17	376	0.38	0.26
$< 45^{\circ}$	Ours-FNL	157	0.15	0.11	367	0.35	0.25	533	0.50	0.36
	Ours-INST	141	0.13	0.10	283	0.27	0.19	388	0.37	0.26
	SubCNN	99	0.10	0.07	257	0.26	0.17	383	0.38	0.26
	Ours-FNL	169	0.16	0.11	425	0.40	0.29	618	0.58	0.42
	Ours-INST	149	0.14	0.10	320	0.30	0.22	450	0.43	0.31
	SubCNN	104	0.10	0.07	272	0.27	0.18	409	0.41	0.28

A chair is detected (a), and later becomes occluded (b, shown in dashed lines). Our system predicts its re-Research sponsored by ARO W911NF-15-1-0564/66731-CS, ONR appearance and resumes update (c). N00014-17-1-2072, AFOSR FA9550-15-1-0229.



Real Scene Results

Videos & Code available at <u>http://vision.ucla.edu/vis</u>